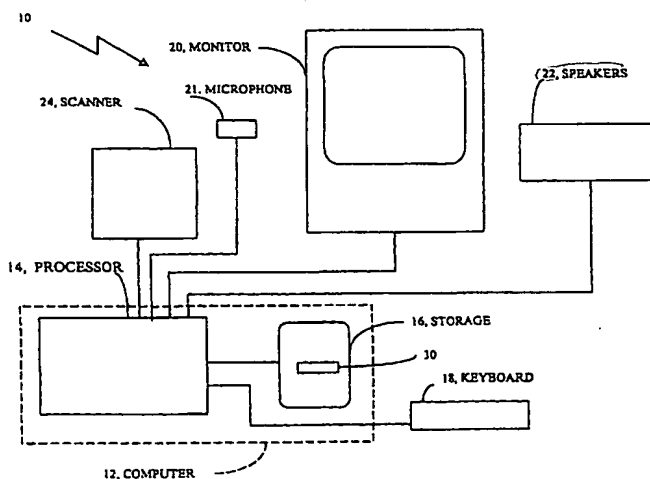




## INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification 6 : <b>G10L 3/00, G09B 5/00</b>	<b>A1</b>	(11) International Publication Number: <b>WO 99/66493</b> (43) International Publication Date: <del>23 December 1999</del> (23.12.99)
<p>(21) International Application Number: PCT/US99/13886</p> <p>(22) International Filing Date: <del>21 June 1999</del> (21.06.99)</p> <p>(30) Priority Data: 09/100,299 19 June 1998 (19.06.98) US</p> <p>(63) Related by Continuation (CON) or Continuation-in-Part (CIP) to Earlier Application US 09/100,299 (CON) Filed on 19 June 1998 (19.06.98)</p> <p>(71) Applicant (for all designated States except US): KURZWEIL EDUCATIONAL SYSTEMS, INC. [US/US]; 411 Waverly Oaks Road, Waltham, MA 02154 (US).</p> <p>(72) Inventor; and (75) Inventor/Applicant (for US only): KURZWEIL, Raymond [US/US]; 203 Lake Avenue, Newton, MA 02161 (US).</p> <p>(74) Agent: MALONEY, Denis, G.; Fish &amp; Richardson P.C., 225 Franklin Street, Boston, MA 02110-2804 (US).</p>	<p>(81) Designated States: AE, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CU, CZ, DE, DK, EE, ES, FI, GB, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, UA, UG, US, UZ, VN, YU, ZA, ZW, ARIPO patent (GH, GM, KE, LS, MW, SD, SL, SZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).</p> <p><b>Published</b> With international search report. Before the expiration of the time limit for amending the claims and to be republished in the event of the receipt of amendments.</p>	

(54) Title: COMPUTER AUDIO READING DEVICE PROVIDING HIGHLIGHTING OF EITHER CHARACTER OR BITMAPPED BASED TEXT IMAGES



## (57) Abstract

A reading system includes a computer (12) and a mass storage device (16) and software including instructions for causing a computer to accept an image file of a document generated from scanner (24). The software converts the image file into a converted text file that includes text information, and positional information associating the text with the position of its representation in the image file. The software records the voice of an operator of the reading machine as a series of voice samples received from microphone (24) in synchronization with a highlighting indicia applied to a displayed representation of the document on a monitor (20) and stores the series of voice samples in a data structure that associates the voice samples with displayed representation. The reading machine plays back the stored, recorded voice samples corresponding to words in the document as displayed by the monitor while highlighting is applied to the words in the displayed document.

**FOR THE PURPOSES OF INFORMATION ONLY**

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece			TR	Turkey
BG	Bulgaria	HU	Hungary	ML	Mali	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MN	Mongolia	UA	Ukraine
BR	Brazil	IL	Israel	MR	Mauritania	UG	Uganda
BY	Belarus	IS	Iceland	MW	Malawi	US	United States of America
CA	Canada	IT	Italy	MX	Mexico	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NE	Niger	VN	Viet Nam
CG	Congo	KE	Kenya	NL	Netherlands	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NO	Norway	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	NZ	New Zealand		
CM	Cameroon			PL	Poland		
CN	China	KR	Republic of Korea	PT	Portugal		
CU	Cuba	KZ	Kazakhstan	RO	Romania		
CZ	Czech Republic	LC	Saint Lucia	RU	Russian Federation		
DE	Germany	LI	Liechtenstein	SD	Sudan		
DK	Denmark	LK	Sri Lanka	SE	Sweden		
EE	Estonia	LR	Liberia	SG	Singapore		

- 1 -

# COMPUTER AUDIO READING DEVICE PROVIDING HIGHLIGHTING OF EITHER CHARACTER OR BITMAPPED BASED TEXT IMAGES

## BACKGROUND

This invention relates to reading machines that  
5 read aloud electronic representations of displayed documents.

Reading machines have been used to improve the educational attainment of individuals with learning disabilities. In general, reading machines are computer-  
10 based having specialized software that processes an input source document and generates synthetic speech. This enables the user to hear the computer read the document a word, line, sentence, etc. at a time. Often these reading machines include a scanner to provide one  
15 technique to input source documents to the reader.

The scanner provides an image file representation of a scanned document. The personal computer using optical character recognition software produces an OCR file including generated text information. The OCR file  
20 is used by the display system software to display a text-based representation of the scanned document on the monitor. The OCR file text is also used by speech ~~synthesis software to synthesize speech. Techniques are known for applying highlighting to displayed text of a~~  
25 ~~document synchronized with synthesized speech corresponding to the highlighted text.~~

## SUMMARY

In one aspect of the invention, a computer program product residing on a computer readable medium includes  
30 instructions for causing a computer to ~~display a representation of a document on a computer monitor.~~ The product also causes the computer to read the displayed representation of the document by using a recorded human voice. Optionally the ~~recorded human voice is~~

- 2 -

~~synchronized with highlighting applied to the displayed representation.~~ The computer program product uses information associated with a text file to synchronize the recorded human voice and the highlighting to the  
5 displayed representation of the document.

With a further aspect of the invention, a computer program residing on a computer readable medium includes instructions for causing a computer to record the voice of an operator of the reading machine as a series of  
10 voice samples in synchronization with a highlighting indicia applied to a displayed document and store the series of voice samples in a data structure in a manner that associates the series of voice samples with displayed positions of words in the document.

15 According to a further aspect, a computer program product residing on a computer readable medium includes instructions for causing a computer to playback recorded voice samples corresponding to words in a document displayed by a monitor of the computer and highlight  
20 words in the displayed document in synchronization with the recorded voice samples.

According to a still further aspect of the invention, a reading system includes a computer which includes a processor and a ~~computer monitor for~~  
25 ~~displaying an image representation of a scanned document.~~

The computer also includes a scanner for scanning documents, ~~speakers for providing an audio output~~ and a mass storage device storing a computer program product including instructions for causing the computer to  
30 display a representation of a document on the computer monitor and apply digitized, voice samples of the document to an audio system to cause the computer to output a human voice pronunciation of the document.

According to a still further aspect, a method of

- 3 -

operating a reading machine includes displaying a representation of a document and using positional information from a text file associated with the document to apply digitized, voice samples of the document to an  
5 audio system causing the reading machine to read the document aloud with a recorded human voice pronunciation of the document.

In this manner, a more pleasing pronunciation of the words in the document is provided. The computer  
10 program product can operate in conjunction with a text or image-based representation of the document.

#### BRIEF DESCRIPTION OF THE DRAWINGS

The foregoing features and other aspects of the invention will be described further in detail by the  
15 accompanying drawings, in which:

FIG. 1 is a block diagram view of a reading machine system;

FIG. 2 is a flow chart showing a process to cause the reading system to display and read aloud from a  
20 scanned image representation of a document;

FIG. 3 is a flow chart showing a process used to associate user-selected text on the displayed image representation with recorded human voice samples;

FIG. 3A-3B are flow charts showing processes for  
25 recording and synchronous playback of a digitized human voice in the reading machine system of FIG. 1;

FIGS. 3C-3D are views of data structures or files used with the processes of FIGS. 3A-3B;

FIGS. 4A-4C are flow charts which show a process  
30 to determine a nearest word for use in the process described in conjunction with FIG. 3;

- 4 -

FIG. 4D is a pictorial illustration of a portion of an image representation of text displayed on a monitor useful in understanding the process of FIGS. 4A-4C;

FIG. 5 is a flow chart showing a process to highlight a selected word for use in the process described in conjunction with FIG. 3;

FIG. 6 is a diagrammatical representation of a data structure used in the process of FIG. 3;

FIGS. 7-9 are diagrammatical views of detailed portions of the data structure of FIG. 6;

FIGS. 10A-10C are flow charts of an alternative embodiment for determining the nearest word; and

FIG. 11 is a pictorial illustration of a portion of an image representation of text displayed on a monitor useful in understanding the process of FIGS. 10A-10C.

#### DETAILED DESCRIPTION

Referring now to FIG. 1, a reading machine 10 includes a computer system 12 such as a personal computer. The computer system 12 includes a central processor unit (not shown) that is part of a processor 14. A preferred implementation of the processor 14 is a Pentium-based system from Intel Corporation, Santa Clara, California, although other processors could alternatively be used. In addition to the CPU, the processor includes main memory, cache memory and bus interface circuits (not shown). The computer system 12 includes a mass storage element 16, here typically the hard drive associated with personal computer systems.

The reading system 10 further includes a standard PC type keyboard 18, a sound card (not shown), a pointing device such as a mouse 19, a ~~monitor 20~~, microphone 21, ~~speakers 22,~~ and a scanner 24 all coupled to various ports of the computer system 10 via appropriate interfaces and software drivers (not shown). The

- 5 -

computer system 12 here operates under a WindowsNT® Microsoft Corporation operating system although other systems could alternatively be used.

Resident on the mass storage element 16 is ~~image~~  
5 ~~display and conversion software 30 (FIG. 2)~~ that controls  
the display of a scanned image provided from scanner 24.  
In addition, the software 30 permits the user to control  
various features of the reading machine 10 by referencing  
the image representation of the document as displayed by  
10 the monitor.

Referring now to FIG. 2, the image display and  
conversion software 30 scans 32 an input document in a  
conventional manner to provide an image file 31. The  
image file 31 is operated on by an optical character  
15 recognition (OCR) module 34. The OCR module 34 uses  
conventional optical character recognition techniques  
(typically software based) on the data provided from the  
scanned image 32 to produce an OCR data structure 35.  
Alternatively, image-like representations can be used as  
20 a source such as a stored bit-mapped version of a  
document.

A preferred arrangement of the output data  
structure is described in conjunction with FIGS. 6-9.  
Suffice it here to say, however, that the array of OCR  
25 data structures generally denoted as 35 includes OCR  
converted text, positional and size information for each  
text element. The positional and size information  
associates each text element to its location in the image  
representation of the document as displayed on the  
30 monitor 20.

Referring momentarily to FIG. 7, it can be seen  
that a data structure element 140 includes for a  
particular word an OCR text representation of the word  
stored in field 142. The data structure 140 also has

- 6 -

positional information including X-axis coordinate information stored in field 143, Y-axis coordinate information stored in field 144, height information stored in field 145 and width information stored in field 146. This positional information defines the bounds of an imaginary rectangle enclosing an area associated with the corresponding word. That is, if a pointer device such as a mouse has coordinates within the area of this rectangle, then the mouse can be said to point to the word within the defined rectangle.

~~The image file 31 is fed to a display system 38 which, in a conventional manner, displays 39 the document represented by the image file on the monitor 20. The text file 35 provides one input along with commands driven by the operating system (not shown) to a module 40 which is used to synchronize highlighting and recorded human speech with an image or text displayed representation of the document.~~ Both the image file 31 and the text file 35 may be stored in the reading system 10 for use during the session and can be permanently stored. The files are stored using generally conventional techniques common to Windows95®, WindowsNT® or other types of operating systems.

Referring now to FIG. 3, the user controls operation of the reading system 10 with reference to the image displayed on the monitor 20 by the software module 40. A user initiates reading of the scanned document at the beginning of the document by selecting a reading mode. Among other options, ~~the user can select to hear the document read aloud using synthesized speech or a recorded digitized human voice. The user can have the reading machine 10 start reading the document from any point in the document by illustratively pointing to the image representation of an item from the scanned document~~



- 7 -

~~displayed 42 on the monitor.~~ The document item is the actual image representation of the scanned document rather than the conventional text file representation. The item can be a single word of text, a line, sentence, paragraph, region and so forth.

In addition to pointing to a word, a pointer such as a mouse can point within the text in an image in other ways that emulate the pointer behavior typically used in computer text displays and word processing programs. For instance, by simply pointing to a word the software selects a position in the text before the word; whereas, pointing to a word and clicking a mouse button twice will cause the word to be selected and pointing to a word and clicking an alternate mouse button selects several words, starting at a previously determined point and ending at the word pointed to.

The user can use a mouse or other type of pointing device to select 42 a particular word. Once selected, the software fetches 44 the coordinates associated with the location pointed to by the mouse 19 (FIG. 1). Using these coordinates, the word or other document item nearest to the coordinates of the mouse is determined. The information in the data structure 100 is used to generate highlighting of the word as it appears on the display item as well as recorded digitized speech samples.

While the user can enable the reading machine to generate synthesized speech, it is preferred that the user enable the reading machine to use a recorded digitized human voice to read the document aloud. For purposes of explanation, it will be assumed that the document item is an image representation of a word that is read aloud using a recorded digitized human voice.

The searching process 46, as will be further

- 8 -

described in conjunction with FIGS. 4A-4C, will search for the nearest word. Alternatively, a searching process 46' as will be described with FIGS. 10A-10C can also be used. The search operation performed by searching  
5 process 46' is based upon various attributes of a scanned image.

After the nearest word or nearest document item has been determined 46 (or 46'), the highlighting is applied 48 to an area associated with the item or word.  
10 The locational information in the text file 35 corresponding to the nearest document item is used to locate in a data structure a stored, recorded digitized human voice pronunciation of the word. The stored samples are retrieved from a data structure 41 or 41' and  
15 are fed to the audio system 54 that reads the word aloud. ~~The word is read aloud as the monitor 20 displays the document and highlighting is applied to the word.~~

Referring now to FIG. 3A, a recording process 150 that records a human voice as continuous voice samples is  
20 shown. The process 150 associates the continuous digitized voice samples with positional information of text in a displayed representation of a document. The recording process 150 displays the document such as an image representation of the document on the display 20  
25 (FIG. 1). An operator of the reading machine 10 will read the document aloud as the reading machine records the operator's voice as a plurality of digitized voice samples which are stored in a file or a data structure 41 (FIG. 3C) or 41' (FIG. 3D). The data structure 41  
30 includes positional information that associates the recorded, plurality of digitized voice samples corresponding to continuous speech to the word in the position of the document as it appears on the display. The data structure 41' includes the plurality of recorded

- 9 -

voice samples and links back to the OCR data structure  
35.

The operator of the reading machine 10 talks into  
a microphone 21 (FIG. 1) in synchronization with a  
5 highlighting indicia applied, a word at a time, to the  
displayed document. The operator attempts to synchronize  
his rate of pronunciation of the word to the rate that  
highlighting is applied to the word. ~~Generally, the~~  
~~highlighting is applied based upon a text to speech~~  
10 ~~conversion process~~, as described in conjunction with FIG.  
5. ~~The highlighting is applied at a rate in accordance~~  
~~with the rate at which the word is pronounced by the~~  
~~synthesized speech.~~ Thus, for example, ~~a word which~~  
~~takes a long time to pronounce will have highlighting~~  
15 ~~applied for a longer period of time than a word which~~  
~~takes a shorter period to pronounce.~~ Also, punctuation  
may affect the rate of highlighting. The accuracy of the  
positional information or links between words and the  
continuous voice samples is related to the extent that  
20 the operator can synchronize his pronunciation to  
highlighting.

The reading machine 10 retrieves 152 the first or  
next word in the document and applies highlighting 154 to  
the word indicating which word the operator should  
25 pronounce into the microphone 21. The pronunciation of  
the word will be recorded and stored as a plurality of  
digitized voice samples by the reading machine 10 using  
standard software drivers for the input of audio signals.  
The digitized voice samples are stored 158 in the data  
30 structure 41 along with positional information from data  
structure 35 or in data structure 41' with links back to  
data structure 35. This positional information or links,  
as appropriate, approximate starting samples of words in  
the document. This process continues at 160 until there

- 10 -

are no more words in the document or until the operator terminates the process. At the end of the process 150 a data structure 41 (FIG. 3C) or 41' (FIG. 3D) is populated with the recorded voice samples.

5           As an alternative to the recording process of FIG. 3A a so-called large vocabulary continuous speech recognition software system available from Kurzweil Applied Intelligence, (a division of Lernout & Hauspie), Dragon Systems, or IBM can be used to transcribe the  
10 recorded speech. These systems, however, would have to be modified so that they would be capable of marking exactly where each word is in the data structures 41 or 41'. The speech recognition is potentially made easier since the transcribed text is provided from the OCR  
15 software. Even if the output of the speech recognition is not perfect, it will identify correctly most of the words. This can be used to improve synchronization links.

Referring now to FIG. 3B, recorded playback of the  
20 voice samples in the reading machine 10 in synchronization with highlighting of words in a displayed document is shown. The process 51 starts with retrieving coordinates of a nearest document item. In general this would be a word pointed to by a mouse. For an image  
25 representation there may be some ambiguity in what is actually being pointed to. Therefore the approaches described below can be used to obtain the nearest document item. The coordinates are used to access the data structures 41 or 41' and extract 154 the appropriate  
30 voice samples. With data structure 41 a match in coordinates is sought, using one of the two alternate methods described below (methods 46 or 46'), whereas with data structure 41', method 46 (described below) is used, by fetching 162 coordinates of successive words by

- 11 -

following each link in data structure 41' to the data structure 35 (FIG. 7) and using the coordinates in data structure 35. The extracted voice samples are sent 166 to a standard software driver for the output of audio signals. The outputted audio signals are fed 168 to the audio system. With the highlighting applied as in 48 (FIG. 3) and the recorded voice samples fed to the audio system, the reading machine 10 reads the document aloud to the user, using the recorded human voice samples, while the reading machine 10 applies highlighting or double highlighting to the displayed representation of the document.

Referring now to FIG. 3C, the data structure 41 includes data structure element 148 in which are stored voice samples associated with the word and positional information such as the X and Y coordinates, and height and width information that associates the word in the document to the stored voice samples. Data structure 41 is similar in design to that of data structure 35. Therefore, data structure 41 includes structures 112, 128 and 138 to represent pages, regions and lines respectively.

As shown in FIG. 3D, data structure 41' stores voice samples stored with links back to the data structure 35. Either the data structure 41 or the data structure 41' associate recorded human voice samples of words in the displayed document to the position of the words in the image or text-based representation of the document.

The positional information in data structure 41 or the links in data structure 41' produced by the above process may be only approximately correct. For many applications this would be sufficient allowing the system

- 12 -

to work adequately most of the time. It is possible, however, to improve on the accuracy of the synchronization between the human voice and the displayed words. One technique would have the user of the system  
5 10 play back the recorded voice samples while displaying the document. The word highlighting (either single or double highlighting) is applied to the words and user speeds up or slows down the highlighting process to better synchronize the highlighting to the playback of  
10 the recorded human voice. For example, the up and down cursor keys or left and right cursor keys or any other pair of keys can be manipulated to accelerate or decelerate the progress of the word highlighting to improve synchronization of the word highlighting with the  
15 recorded human. The user can press the up arrow key and while keeping the key depressed accelerate highlighting. Release of the key would cause the word highlighting to return to its original rate.

Alternatively, the speech rate can be changed.  
20 The literature on audio signal processing describes several methods of time-scale modification which can be used to allow the recorded voice samples to be played back at a slower or faster rate without substantially changing the pitch of the voice.

25 Alternatively, the reading machine 10 can include a visual sound editor that permits the user to play a segment of the voice recording. The operator can then identify the word or words corresponding to the voice recording and correct the positional information or the  
30 link, as appropriate, for that word or words in the data structures 41, 41'.

Optionally, standard data compression and decompression techniques can be used to store the voice samples.

- 13 -

Described below are processes used to determine a nearest word in the image as well as a process used to highlight a word or apply double highlighting to a word. In essence, these processes can operate on a display of the document by use of the image file. The software makes reference to the OCR data structure 35 to determine positional information to associate the reading software, highlighting software or other software with respect to commands by the user. The above data structures 41 or 10 41' can be saved in a file for later use.

Referring now to FIGS. 4A-4C, the process 46 used to determine a nearest word in an image display as pointed to by a user is shown. A pointer is initialized 60 and a maximum value is loaded into a displacement 15 field 51b of structure 51 (FIG. 4D). The displacement field 51b is used to store the smallest displacement between a word boundary and the coordinates of the pointing device. The pointer that is initialized 60 is a pointer or index into the OCR generated data structure 35 20 (FIG. 6). The software 46 retrieves each word entry in the data structure 35 to determine for that word, in accordance with the image, relative position information associated with the OCR text generated word whether or not that particular word is the closest word to the 25 coordinates associated with the user's pointing device.

The coordinates associated with a first one of the words are fetched 62. The coordinates associated with the first one of the fetched words are used to determine 64 whether the pointing device is pointing to a location 30 within a box 65, that is defined around the word. Thus, as shown in conjunction with FIG. 4D, the mouse points to a spot 61 having coordinates  $X_i$ ,  $Y_j$ . For any document item on the scanned image, an imaginary box here 65, is assumed to exist about the word "IMAGE" in FIG. 4D.

- 14 -

Thus, if the pointing device coordinates fall within the box 65<sub>5</sub>, the pointing device would be considered to point to the document item "IMAGE" associated with the box 65<sub>5</sub>.

In the data structure 35 each of the words will have associated therewith the OCR text converted from the image file 31, as well as position and size data that identifies the position and size of the word as it appears on the original document. Accordingly, this information also locates the word as it appears in the displayed image representation of the document. Thus, when determining the closest word to a position pointed to by a mouse, it is necessary to determine the boundaries of the box occupied by the particular word. The software determines 64 whether or not point 61 falls within the box by considering the following:

For a mouse coordinate position (X, Y) the location pointed to by the mouse can be considered to be within a region of an image word having points defined by coordinates (a<sub>i</sub>, b<sub>j</sub>) and (c<sub>k</sub>, d<sub>l</sub>) where  $c_k = a_i + h$  and  $d_l = b_j - w$ , if  $X \geq a_i$  and  $Y \leq b_j$  and  $X \leq c_k$  and  $Y \geq d_l$  where it is assumed here that the positive direction of the coordinates is upward and to the right.

If this condition is satisfied, then the point 61 can be considered to be within the box and, hence, control will pass 66 directly to 50 (FIG. 4B). From the information mentioned above, therefore, the point (c, d) can be determined by subtracting the height of the box from the x coordinate (a<sub>i</sub>) associated with the image and adding the width of the box associated with the y coordinate (b<sub>j</sub>) of the image.

If, however, the point 61 is not within the box as is shown, then the software 46 determines 68 the word which is nearest to the point 61 by one of several algorithms. A first algorithm which can be used is to



- 15 -

compute the distance from a consistent corner of the box associated with the word to the position of the mouse pointer 61. In general, the distance (S) to a consistent corner would be computed as the "Pythagorean" technique  
5 as follows:

$$S = ((X-a_i)^2 + (Y-b_j)^2)^{-2}$$

Alternatively, this equation can be used at each corner of each box and further processing can be used to determine which one of the four values provided from each  
10 corner is in fact the lowest value for each box.

In either event, the computed value (S) is compared to the previous value stored in displacement field 51b. Initially, field 51b has a maximum value stored therein and the smaller of the two values is  
15 stored 72 in field 51b. Accordingly, the first computed value and the index associated with the word are stored in the structure 51, as shown in FIG. 4C. It is determined 74 whether or not this is the end of the data structure 35. If it is the end of the data structure 35  
20 then control branches to 50 and hence to 52. If it is not the end of the data structure 35 then the pointer is incremented 76 and the next word in the data structure as determined by the new pointer value is fetched 62.

The second time through the process 46 in general  
25 will be the same as the first time except that the process 46 will determine 72 whether the previously stored value ( $S_p$ ) in fields 51a, 51b is greater than or less than a current calculated value ( $S_c$ ) for the current word. If the current value ( $S_c$ ) is less than the previous  
30 value  $S_p$ , then the current value replaces the previous value in field 51b and the index associated with the current value replaces the previous index stored in field 51a.

In this manner, the structure 51 keeps track of

- 16 -

the smallest calculated distance (S) and the index (i.e., word) associated with the calculated distance. The process continues until the positional data for all of the words in the data structure associated with the particular image have been examined. The values which remain in the data structure 51 at the end of the process correspond to the closest word to the location given by the mouse coordinates 61.

Referring now back to FIG. 3, once the nearest coordinates for the nearest data item are determined, the process 40 applies highlighting as appropriate to the selected item. One technique for providing highlighting would simply highlight a line or a paragraph in the text representation displayed on the monitor. The highlighting would be of the current word that is being read aloud to the user. Although this is acceptable, a preferred approach as described herein applies double highlighting and still preferably applies double highlighting to an image representation of a scanned document.

~~The selected paragraph or sentence is highlighted with a first transparent color. Each individual word as read aloud by the recorded digitized voice samples is highlighted with a second, different transparent color.~~  
Accordingly, highlighting is applied 48 in a manner as will now be described.

Referring now to FIG. 5, the highlighting process 48 includes waiting 80 for an event by the software 48. The event is typically an operating system interrupt-type driven operation that indicates any one of a number of operations such as a user of the ~~reading machine 10~~ ~~initiating speech synthesis of a word, sentence or paragraph.~~ The highlighting process 48 remains in that state until an event occurs. When an event occurs all

- 17 -

previous highlighting is turned off 82. The previous highlighting is turned off by sending a message (not shown) to the display system 38 causing the display system to remove the highlighting. The highlighting process checks 84 whether a unit of text has been completed. For example, a unit can be a word, line, sentence, or a paragraph, for example, as selected by the user.

If a unit of text has been completed, then highlighting of the unit is also turned off 90. The software checks for an exit condition 91 after the coordinates have been fetched. An exit condition can be any one of a number of occurrences such as reaching the last word in the array of OCR data structures 35 or a user command to stop. If an exit condition 92 has occurred, the routine 48 exits to 92.

If an exit condition has not occurred, the next unit is determined 93. The next unit of text is determined by using standard parsing techniques on the array of OCR text structures 35. Thus, the next unit is determined by looking for periods, for example, to demarcate the end of sentences, and indents and blank lines to look for paragraphs. In addition, changes in the Y coordinate can be used to give hints about sentences and lines. Other document structure features can also be used. The next unit is then highlighted 94 by instructing the display system software 38 (FIG. 2) to apply a transparent color to the selected next unit. This is a first level of highlighting provided on a unit of the image representation of the scanned document. Control transfers back to 86.

The coordinates of the next word are fetched 86. The software checks 88 for an exit condition after the coordinates have been fetched. An exit condition can be

- 18 -

any one of a number of occurrences such as reaching the last word in the array of OCR data structures 35 or a user command to stop provided from the keyboard 18 or other input device. If an exit condition has occurred  
5 88, the routine 48 exits 89. Otherwise, a second highlight is applied 96 to the image, here preferably with a different transparent color and applied only to the word that is read aloud in the recorded digitized human voice. The pointer to the next word in the data  
10 structure 35 is incremented 98 to obtain the next word. The second highlighting is provided by sending a message to display system software 38 containing the positional information retrieved from the data structure. This process continues until an exit condition occurs 88.

15 It should be noted that the single and the dual highlighting above were described as applying two distinct, transparent colors to the image representation of the displayed document. Alternatively, however, ~~other highlighting indicia can be used such as bold text, font style or size changes, italics, boxing in selected text, and underlining. In addition, combinations of these other indicia with or without colors could be used.~~

Referring now to FIGS. 6-9, a preferred format for the data structure 35, as provided by the OCR 34, is  
25 shown. The data structure 35 is hierarchically organized. At the top of the data structure is a page, data structure 110. The page includes pointers 110a-110e to each one of a plurality of regions 120. A region is here a rectangular-shaped area comprising one or more  
30 rectangular lines of text. If there are multiple lines of text in a region, the lines do not overlap in the vertical direction. That is, starting with the top line, the bottom of each line is above the top of the next line. Here the regions may include headers, titles,

- 19 -

columns and so forth. The headers may or may not straddle more than one column and so forth. The regions likewise include a plurality of pointers 120a-120e to each one of corresponding lines 130 shown in the data structure 130. The lines correspondingly have pointers 130a-130e to each of the words contained within the line.

As shown in conjunction with FIGS. 7-9, the detail structure of items 140, 130 and 120 include a plurality of fields. Thus, for example, FIG. 7 for the word includes the text field 142 which has the OCR generated text and has fields 143 and 144 which provide rectangular coordinate information x and y, respectively, as well as fields 145 and 146 which provide here height and width information. Similar data are provided for the lines as shown in FIG. 8 as well as regions as shown in FIG. 9.

Now to be described will be a preferred method 46' to determine the nearest word associated with the position of a mouse or other pointing device. This approach is particularly advantageous for those situations where dragging operations of a mouse are often performed. The image representation may not provide an exact correspondence to the text as determined by the OCR recognition system. Also sometimes incorrect text is selected because the user does not precisely place the mouse or other pointing device directly on the desired item in the image representation. Also, when the pointer is positioned in the white space between lines, or in the white space to the left or right of lines, choosing the closest word to the pointer will not always give the result that a computer user would normally expect, based on the behavior of mouse selection on standard computer text displays. Moreover, minor misalignments may also occur between the image representation as displayed on the display and as provided by the OCR text file.

- 20 -

Thus, for example, consider point 61c on Figure 11. In the method 46 previously described, the closest word, which is "OF" in the previous line, will be chosen as the selected word. But on standard computer displays  
5 the point of selection would be after the word "LAST."

The approach as shown in conjunction with FIGS. 10A-10C will tend to mitigate some of these errors.

Referring now to FIG. 10A, pointers are again initialized 180 to a first one of the regions and the  
10 coordinates of the region's boundary box are fetched from the data structure 120. The position (X, Y) of the pointer is calculated to determine whether or not it falls within a box defining a region.

To further illustrate this process, reference is  
15 also made to FIG. 11 which shows a sample region containing a plurality of lines of text in the image-based representation and boxes illustrated about the region, lines and word. Also three sample positions 61, 61a, 61b of the pointing device (not shown) are  
20 illustrated.

The calculation for a region is performed in a similar manner as for calculating a box for a word described in conjunction with FIGs. 5A to 5C except that the positional information contained within the region  
25 data structure 120 is used to determine a box or other boundary associated with the region. Coordinates  $(r_s, s_s)$  and  $(t_s, u_s)$  denote the imaginary box about the illustrated region in FIG. 11. If it is determined 186 that the coordinates of the pointer fall within the box (as 61 and  
30 61a-61d, FIG. 11), then the process branches to determine 201 (FIG. 10B) the nearest line 201 (FIG. 10B). Otherwise processing continues to determine 187 whether or not the process has reached the last region in the region data structure 120. If it has not reached the

- 21 -

last structure, the pointer is incremented 194 to point to the next region in the data structure 120. If the process 46' has reached the last structure the coordinates of the pointer device do not point to any word, as 61, (FIG. 11). Therefore, a previously determined word is used, and the process exits.

If it was determined 186 that the coordinates fall within a region's box, then a similar process 201 is used to determine the nearest line except that the line data associated with the data structure 130 (FIG. 8) is used for positional information and index information such as coordinates  $(l_4, m_4)$  and  $(n_4, o_4)$ . Again for each line within the particular region, positional information is used to determine whether the coordinates of the pointing device are within a box defined about the line by the positional information associated with the line. If the coordinates of the positioning device fall above the box associated with the line as point 61a, then the software will choose the first word of the line, here the word "TEXT." If the coordinates fall above the bottom of the line box as point 61b, then the software branches to 220.

As shown in conjunction with FIG. 10B, the software initializes 201 a pointer to the top line in the region and fetches 202 the coordinates of the line. The coordinates which are fetched correspond to the top and bottom coordinates of an imaginary box positioned about the line. The software calculates 204 whether the Y coordinate of the pointing device is above the line. This is accomplished by comparing the value of the Y coordinate of the pointing device to the Y coordinate  $(m_4)$  of the uppermost point defining the box about the line, as shown for point 61b. If it is determined 206 that the Y coordinate is above the box defined about the line, the software chooses 208 the first word on the line and is

- 22 -

done. Otherwise, the software determines whether the Y coordinate is above the bottom of the box defining the line by using a similar approach as for the top of the line except using, for example, the coordinate  $(0_4)$ . If  
5 it is determined that the Y coordinate is equal to or above the bottom of the box defining the line, as point 61b then the software branches to 220 (FIG. 10C).

The X coordinate of the pointer is already known to be in the region and is not checked here. This allows  
10 for short lines to be detected. Lines are often shorter than the width of the region. For example, short lines may occur at the beginning and end of paragraphs or in text that is not justified to form a straight right margin. Otherwise, it continues to 212 where it is  
15 determined whether the current line is the last line in the data structure 230. If it is not the last line in data structure 230, the pointer is incremented 216 to point to the next lower line in the region. If it is the last line in the data structure and the Y coordinate was  
20 not above the top of the line nor above the bottom of the line, the software chooses 214 the word after the word in the last line as for point 61c and is done.

Referring now to FIG. 10C, pointers are again initialized to a first one of the words on a line, as  
25 shown by 220 and the coordinates of the word are fetched 222 from the data structure 140. The position X of the pointer is calculated at 224. This calculation is performed to determine whether or not the X position of the pointer falls at or to the left of the current word's  
30 right side, as shown for point 61a. This calculation is performed by comparing the X value of the pointer coordinate to the X value of the right side of the box defined about the word here coordinate  $a_5$  of point  $(a_5, b_5)$ . If the value of the X coordinate for the box is



- 23 -

less than or equal to that of the X coordinate of the pointing device, then the pointing device is considered pointing to the left side of the word's right side. It is determined 226 whether the pointer points to the left side of the word's right side. If it does, the particular word "TEXT" is chosen 227 for point 61d and the process is done. Otherwise, the process determines 228 whether or not it has reached the last word in the data structure 140. If it has not reached the last word in the data structure 140, the pointer is incremented 234 to point to the next word to the right. If it has reached the last word in the data structure 140, the software will choose 230 the word after the last word in the line (not illustrated) and the process is done.

15       The chosen word is forwarded to 48 of FIG. 3. In this manner double highlighting, as described in conjunction with FIG. 5, is performed on the word chosen by this process. The reading machine can read the word aloud using synthesized speech or a recorded human voice, as also described above.

#### Other Embodiments

It is to be understood that while the invention has been described in conjunction with the detailed description thereof, the foregoing description is intended to illustrate and not limit the scope of the invention, which is defined by the scope of the appended claims. Other aspects, advantages, and modifications are within the scope of the following claims.

What is claimed is:

- 24 -

# CLAIMS

1. A computer program residing on a computer readable medium comprising instructions for causing a computer to:
  - display a representation of a document on a
  - 5 computer monitor;
  - use information from a text file associated with the document to apply digitized, recorded voice samples of the document to an audio system to cause the computer to output a human voice pronunciation of the document.
- 10 2. The computer program product of claim 1 wherein the information is positional information associated with the text file.
3. The computer program product of claim 1 wherein the information is links back to the text file.
- 15 4. The computer program product of claim 1 wherein the computer causes a text-based representation of the document to be displayed.
5. The computer program product of claim 1 wherein the computer causes an image representation of the
- 20 document to be displayed.
6. The computer program product of claim 5 wherein the text file is derived from an image file used to display the image representation of the document.
7. A computer program residing on a computer readable
- 25 medium comprising instructions for causing a computer to:
  - accept an image file generated from optically

- 25 -

scanning a document;

convert an image file into a converted text file,  
said converted text file including text information and  
positional information associating the text with the

5 position of its image file representation;

display the image representation of the scanned  
document on a computer monitor;

select a document item from the displayed image  
representation of the document by using the positional  
10 information in the converted text file; and

read the document item aloud using digitized,  
voice samples of the document item.

8. The computer program as recited in claim 7 wherein  
said program further comprises instructions for causing  
15 the computer to:

feed electrical signals an audio system to cause  
the computer to output the digitized, voice samples.

9. The computer program as recited in claim 7 further  
comprising computer instructions for causing the  
20 displayed image representation of the document item to  
be:

highlighted by applying a highlighting indicia to  
the displayed image representation in accordance with  
positional information provided from the converted text  
25 file.

10. The computer program as recited in claim 7 further  
comprising computer instructions for causing the computer  
to:

highlight the displayed image representation of  
30 the document item with a color by applying a color to the  
displayed image representation in accordance with

- 26 -

positional information provided from the converted text file.

11. The computer program of claim 7 further comprising computer instructions for causing the computer  
5 to:

apply a first highlight to a portion of the document and a second different highlight to each word in the portion of the document as the word is read aloud.

12. A computer program residing on a computer readable  
10 medium comprising instructions for causing a computer to:

record the voice of an operator of the reading machine as a plurality of voice samples in synchronization with a highlighting indicia applied to a displayed document; and

15 store the plurality of voice samples in a data structure in a manner that associates the plurality of voice samples with displayed positions of words in the document.

13. The computer program product of claim 12 further  
20 comprising instructions for causing a computer to:

apply highlighting to the displayed document a word at a time to indicate the word the operator should pronounce into a microphone.

14. The computer program product of claim 12 wherein  
25 the voice samples are stored along with positional information obtained from a text-based data structure.

15. The computer program product of claim 14 wherein the voice samples are stored along with links back to a text-based data structure that contains the positional

- 27 -

information of the recorded words.

16. The computer program product of claim 14 further comprising instructions for causing the computer to playback voice samples, and use a visual sound editor to  
5 permit an operator to identify a word or words corresponding to the voice samples and correct the positional information or links back to the text-based data structure for the word or words.

17. The computer program product of claim 13 wherein  
10 the record process uses a large vocabulary continuous speech recognition software system to transcribe the recorded speech, and mark where voice samples for words are in a voice sample data structure.

18. The computer program product of claim 13 wherein  
15 the speech recognition software uses transcribed text provided from applying optical character recognition to an image representation of a document to produce the text file.

19. A computer program product residing on a computer  
20 readable medium comprising instructions for causing a computer to

playback recorded voice samples corresponding to words in a document displayed by a monitor of the computer; and

25 highlight words in the displayed document in synchronization with the recorded voice samples.

20. The computer program product of claim 19 wherein playback includes instructions to cause the computer to:  
retrieve coordinates of a nearest document item

- 28 -

pointed to by a pointing device; and

retrieve voice samples in accordance with the  
coordinates of the nearest document item.

21. The computer program product of claim 20 wherein  
5 the playback includes instructions to cause the computer  
to:

produce electrical signals in accordance with the  
voice samples; and

feed the electrical signals to an audio system to  
10 read the document using the voice samples.

22. The computer program product of claim 21 wherein  
voice samples are associated with positional information  
in a text-based data structure and further includes  
instructions for causing the computer to:

15 change the speed of highlighting to better  
synchronize the highlighting to playback of the voice  
samples.

23. The computer program product of claim 22 wherein a  
pair of keys on a computer keyboard are manipulated to  
20 accelerate or decelerate a rate of highlighting.

24. The computer program product of claim 23, further  
comprising instructions for causing a computer to:

adjust the rate at which the recorded voice  
samples are played back at a slower or faster rate  
25 without substantially changing the pitch of the speech.

25. A computer program residing on a computer readable  
medium comprising instructions for causing a computer to:  
scan a document to form an image file;  
display on a computer monitor of a computer system

- 29 -

a representation of the document by using the image file;

convert the image file to a text file having text information and positional information; and

in response to a user using a positioning device,

5 cause the computer to select a portion of the displayed image of the document;

retrieve coordinate information corresponding to the position of a cursor provided as a result of the pointing of the pointing device;

10 search the retrieved coordinates to find coordinates corresponding to a nearest item in the image representation of the document;

apply highlighting to the nearest item;

15 extract the text corresponding to the nearest item; and

using positional information from the text file apply digitized, recorded voice samples of the document to an audio system to cause the computer to output a human voice pronunciation of the document.

20 26. The computer program as recited in conjunction with claim 25 wherein instructions to apply the recorded voice samples include instructions to cause the computer to:

25 retrieve coordinates of a nearest document item pointed to by a pointing device; and

retrieve the plurality of voice samples corresponding to the nearest document item as determined by the retrieved coordinates.

27. The computer program product of claim 26 wherein  
30 the instructions to apply recorded voice samples includes instructions to cause the computer to:

produce electrical signals in accordance with the

- 30 -

extracted voice samples; and

feed the electrical signals to an audio system to read the document using a recorded human voice aloud to the user.

5 28. The computer program product of claim 27 wherein voice samples are associated with positional information in the text-based data structure and further includes instructions for causing the computer to:

change the speed of highlighting to better  
10 synchronize the highlighting to playback of the recorded human voice.

29. The computer program product of claim 28 wherein a pair of keys on a computer keyboard are manipulated to accelerate or decelerate the rate of highlighting.

15 30. A reading system comprising:  
a computer, said computer comprising:  
a processor;  
a computer monitor for displaying a  
representation of a document;  
20 speakers for providing an audio output; and  
a mass storage device, said storage device  
including:

software comprising instructions for causing  
the computer to display the representation of the  
25 document on the computer monitor; and  
apply digitized, voice samples of the  
document to an audio system to cause the computer  
to output a human voice pronunciation of the  
document.



- 31 -

31. The reading system of claim 30 wherein the instructions that cause the computer to apply the voice samples use positional information from a text file associated with the document.

5 32. The reading system of claim 31 wherein the computer program further comprises instructions for causing the computer to:

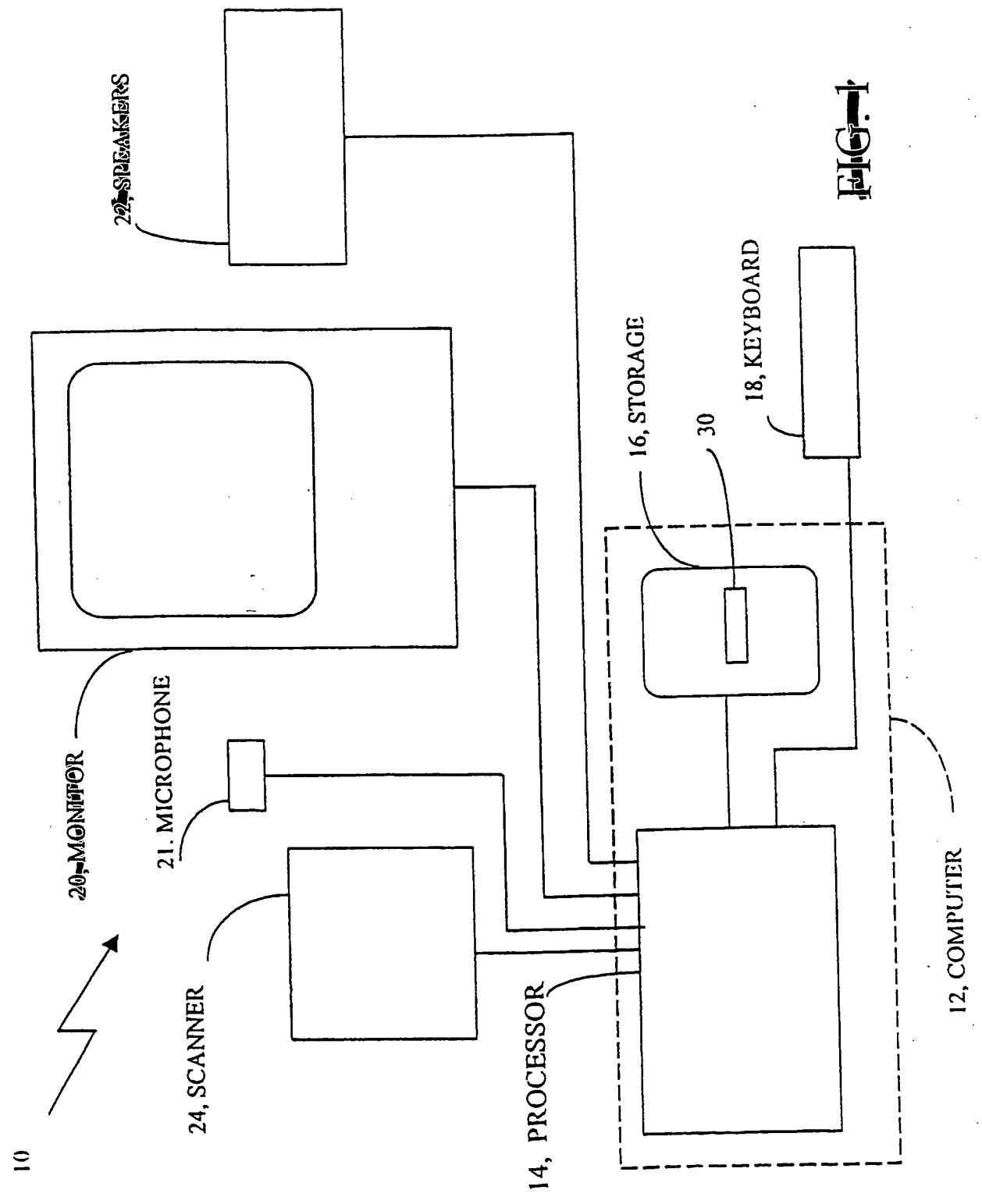
highlight the document in synchronization with the recorded voice samples by applying a highlighting indicia  
10 to the displayed image representation in accordance with positional information provided from the converted text file.

33. A method of operating a reading machine comprises:  
displaying a representation of a document; and  
15 using positional information from a text file associated with the document to apply digitized, voice samples of the document to an audio system to cause the reading machine to read the document aloud with a recorded human voice pronunciation of the document.

20 34. The method of claim 33 wherein the reading machine displays a text-based representation of the document.

35. The computer program product of claim 33 wherein the reading machine displays an image representation of the document.

25 36. The computer program product of claim 34 wherein the text file is derived from an image file used to display the image representation of the document.



**FIG. 1**

2/17

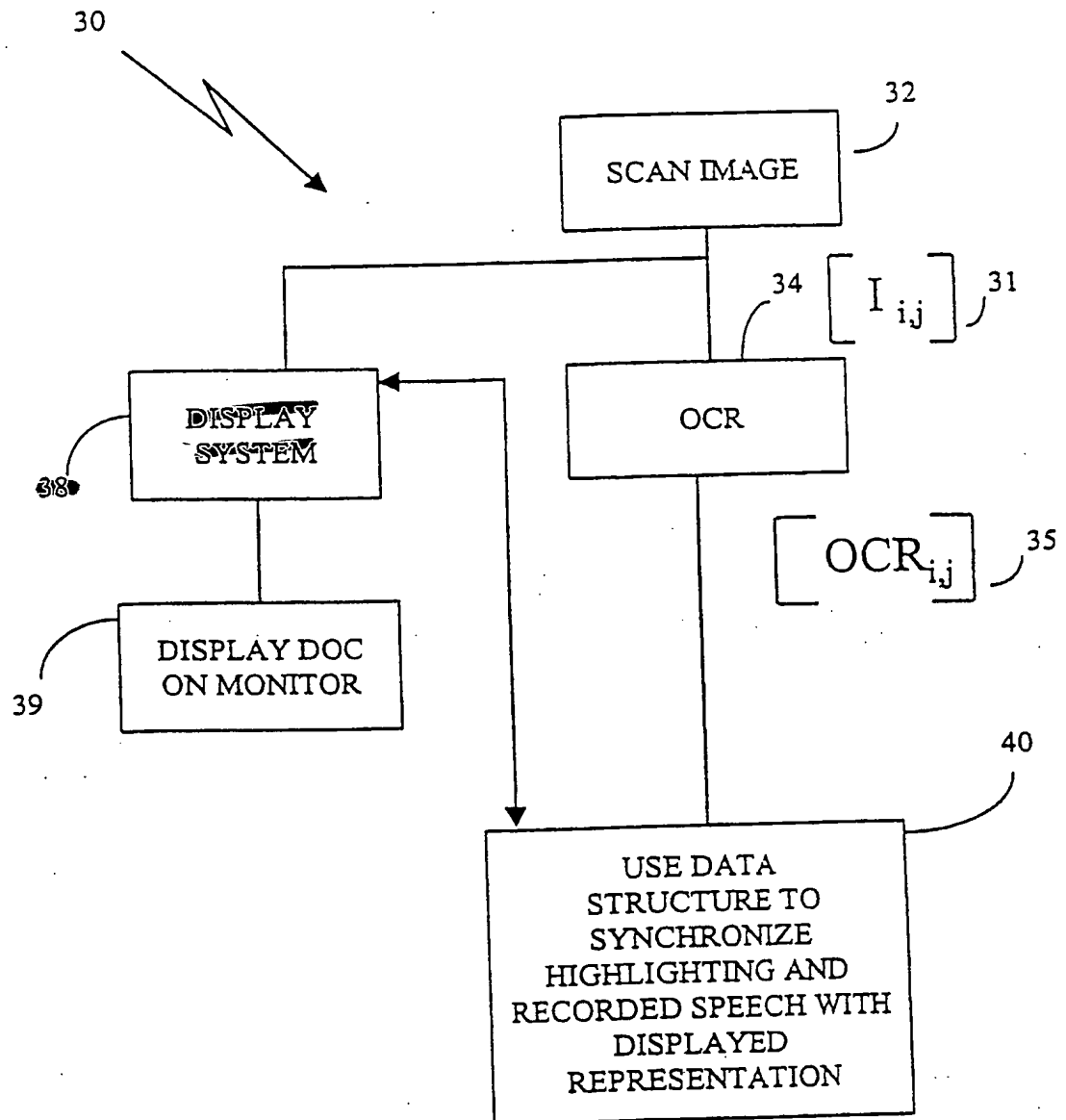
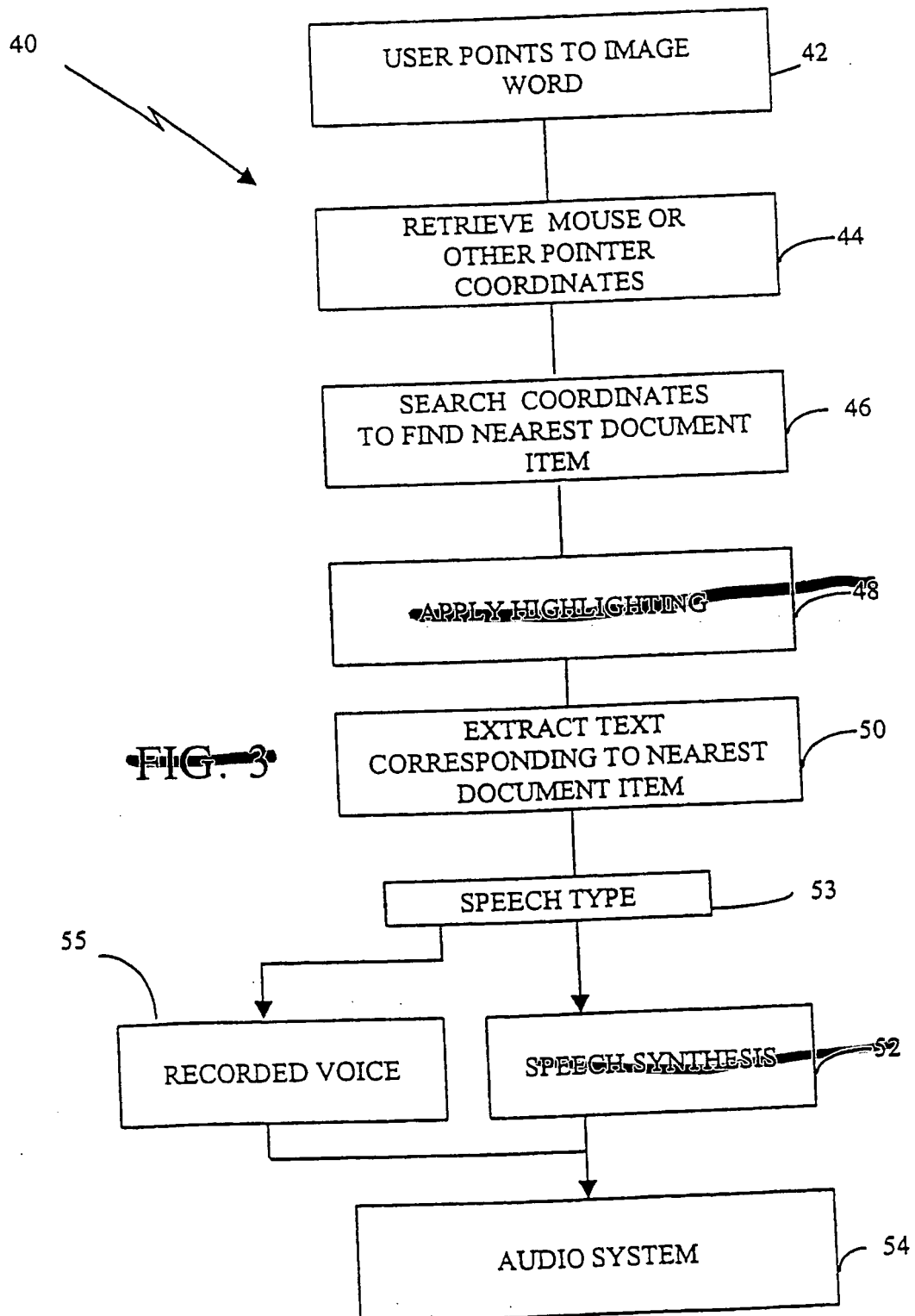


FIG. 2

3/17



4/17

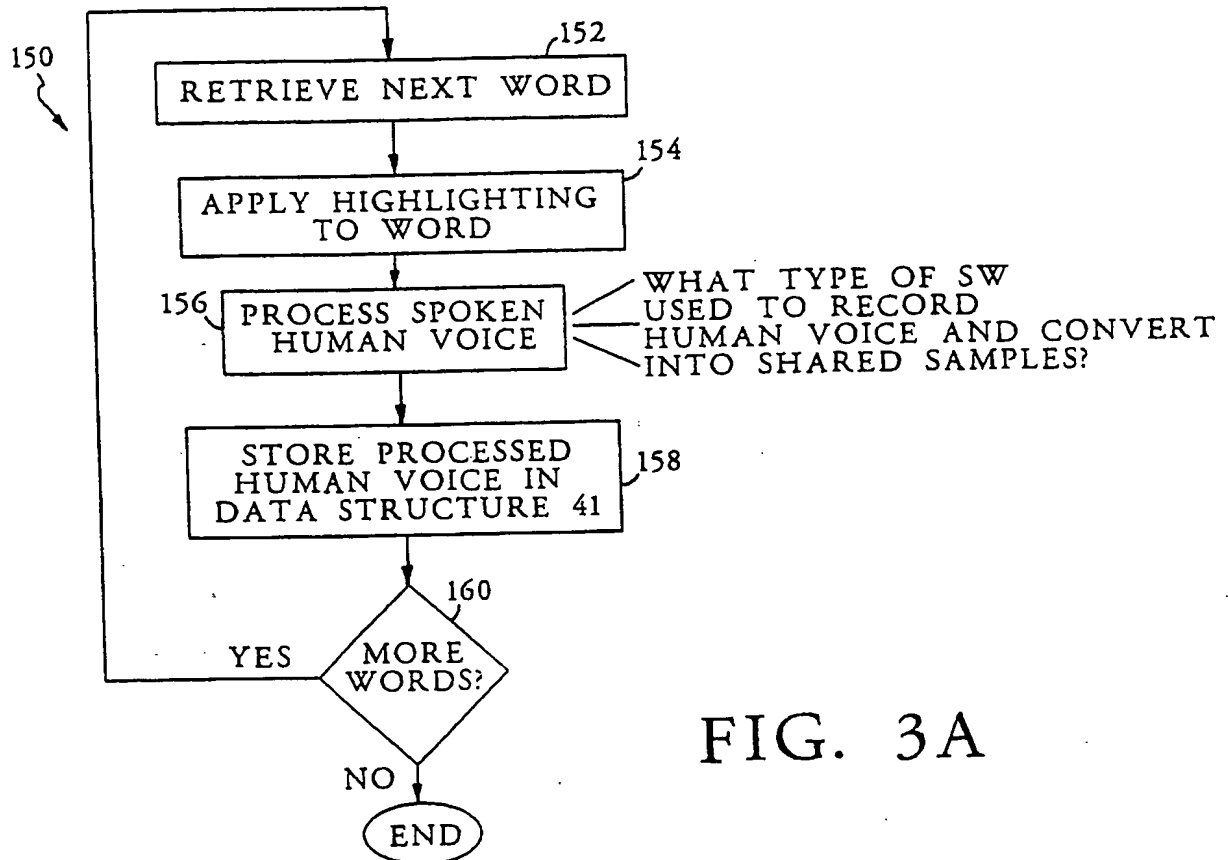
RECORD PROCESS

FIG. 3A

5/17

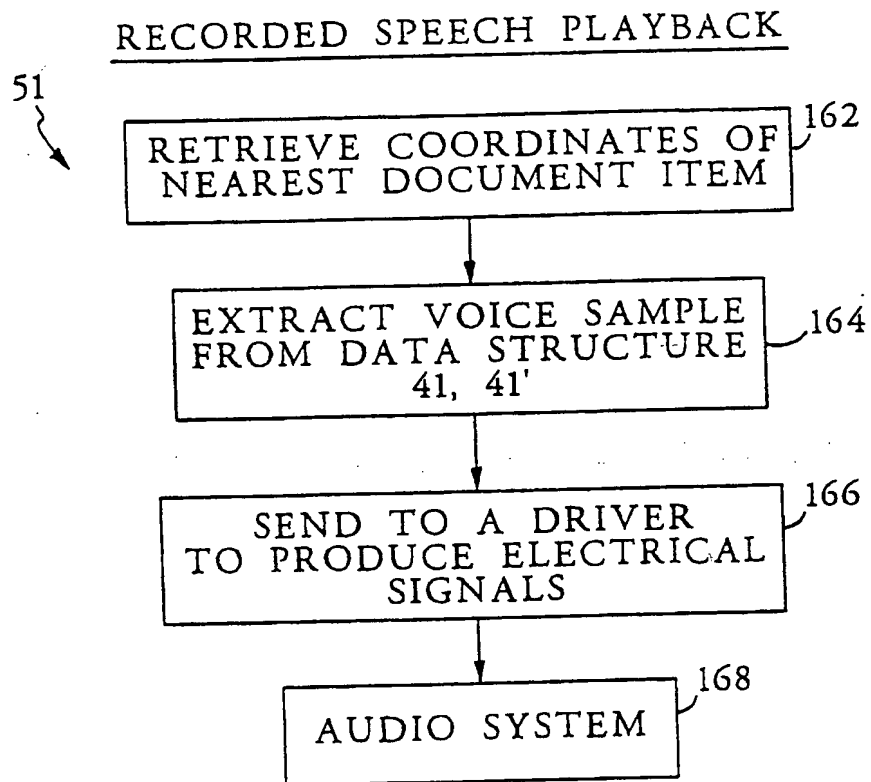


FIG. 3B

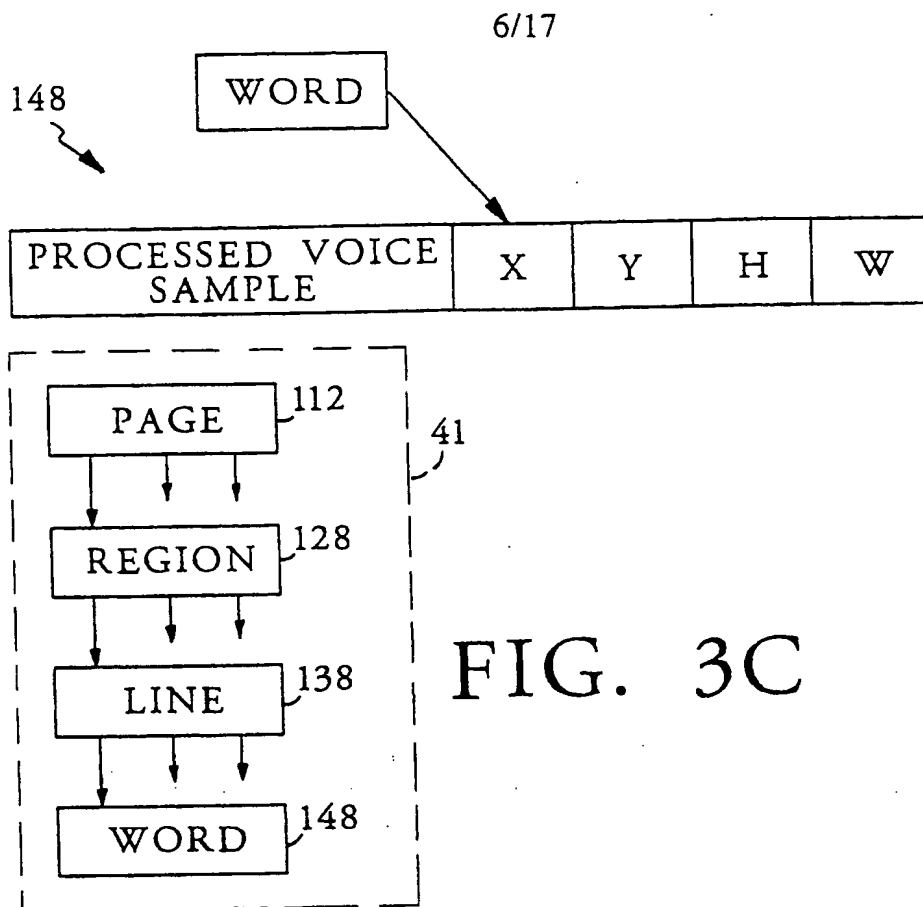


FIG. 3C

RECORDER SAMPLE	DATA STRUCTURE 35
SAMPLE i	ENTRY IN 35

41'

FIG. 3D

7/17

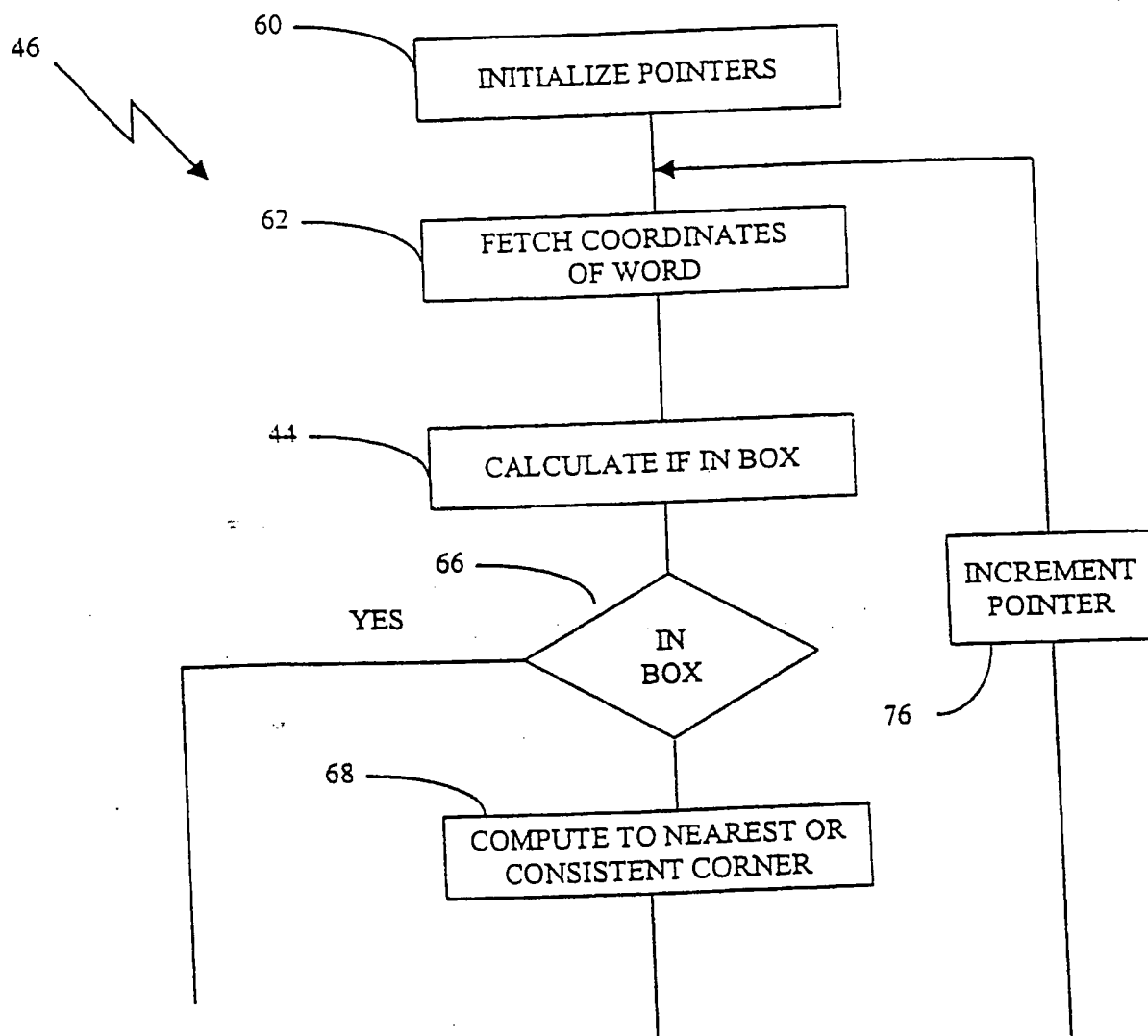


FIG. 4A



8/17

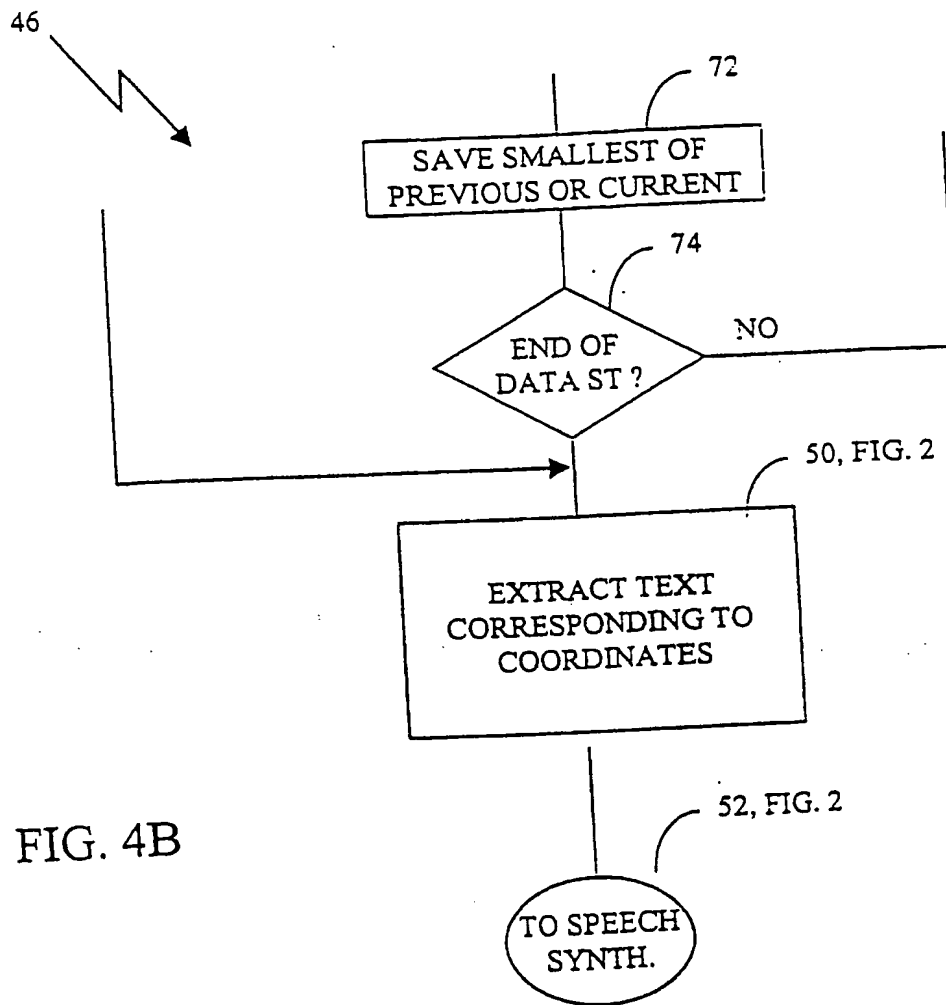


FIG. 4B

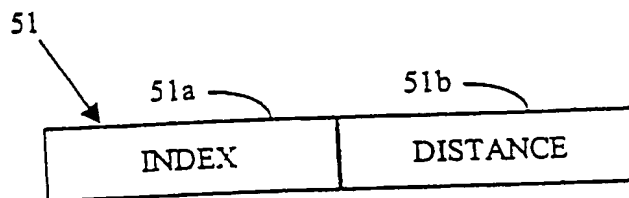
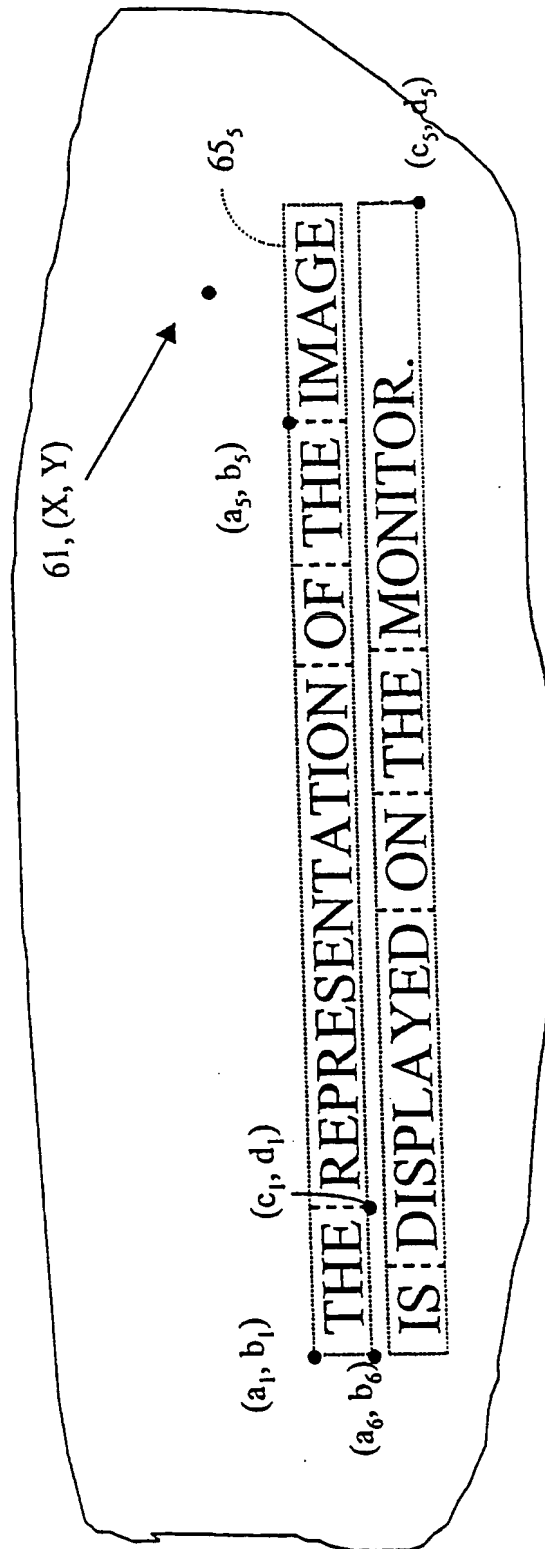
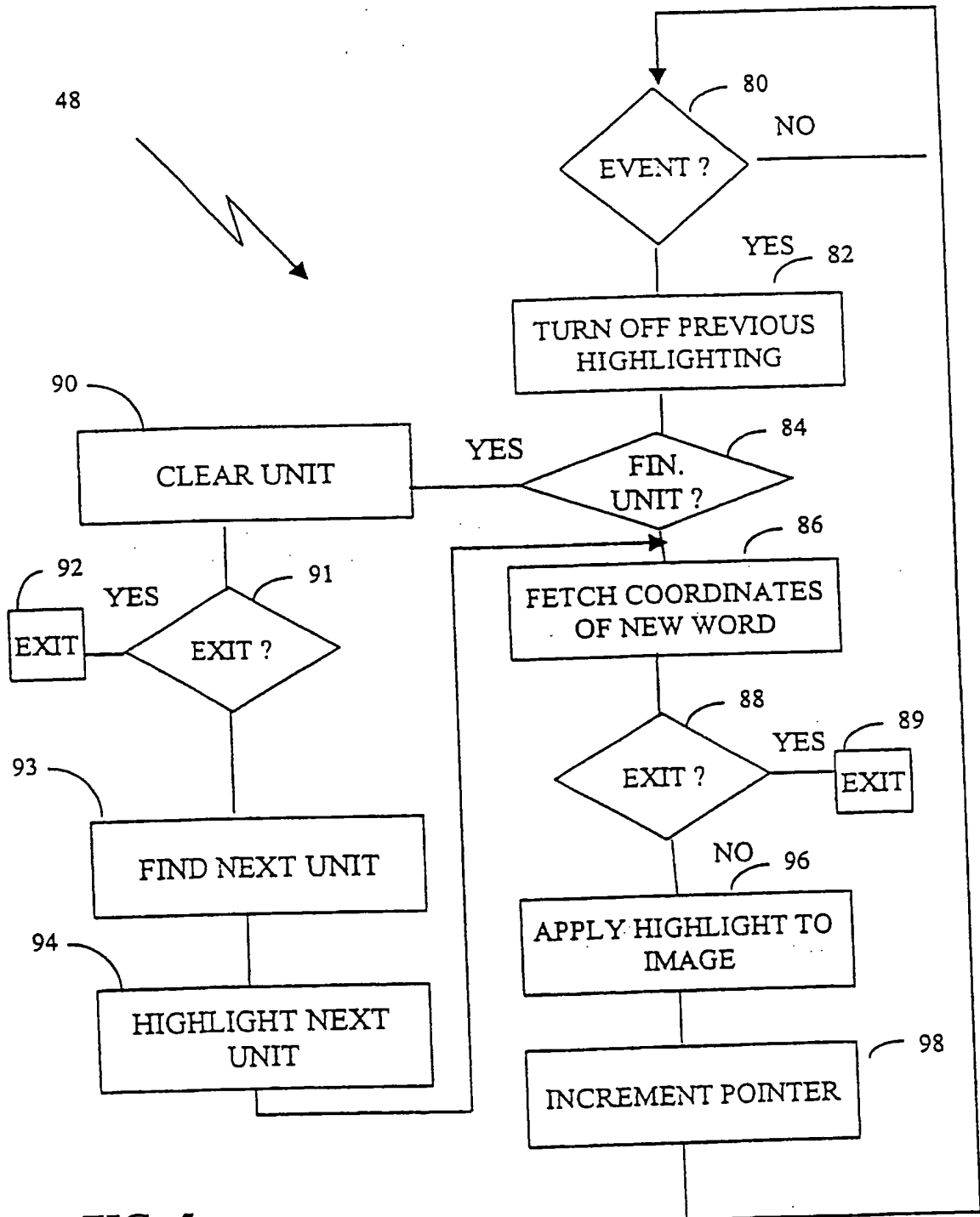


FIG. 4C

9/17

**FIG. 4D**

10/17



11/17

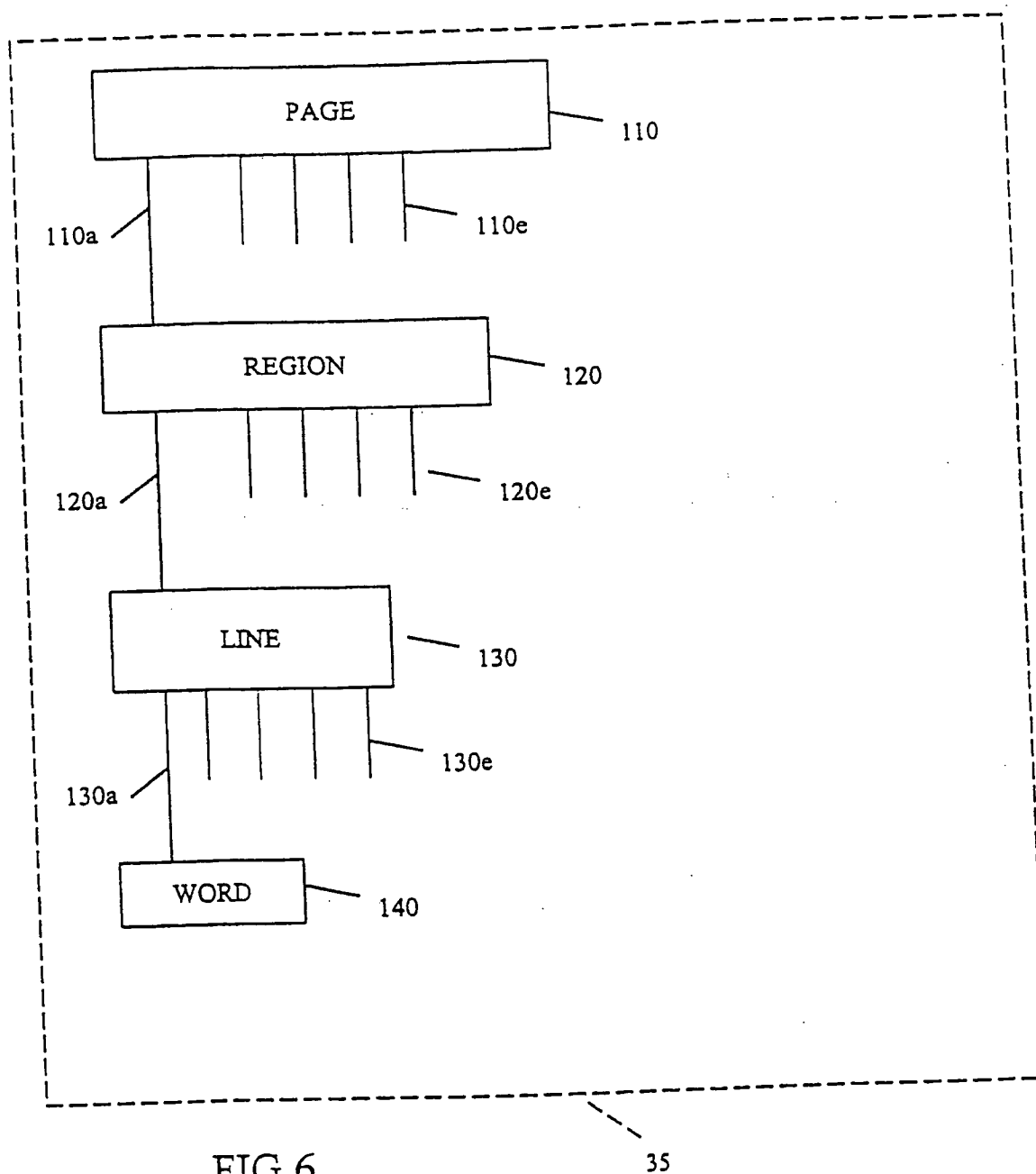


FIG.6

35

12/17

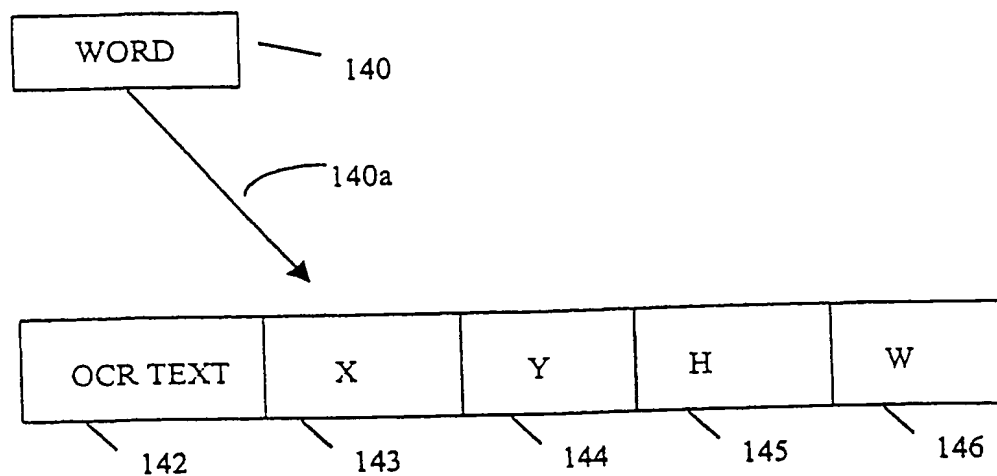


FIG. 7

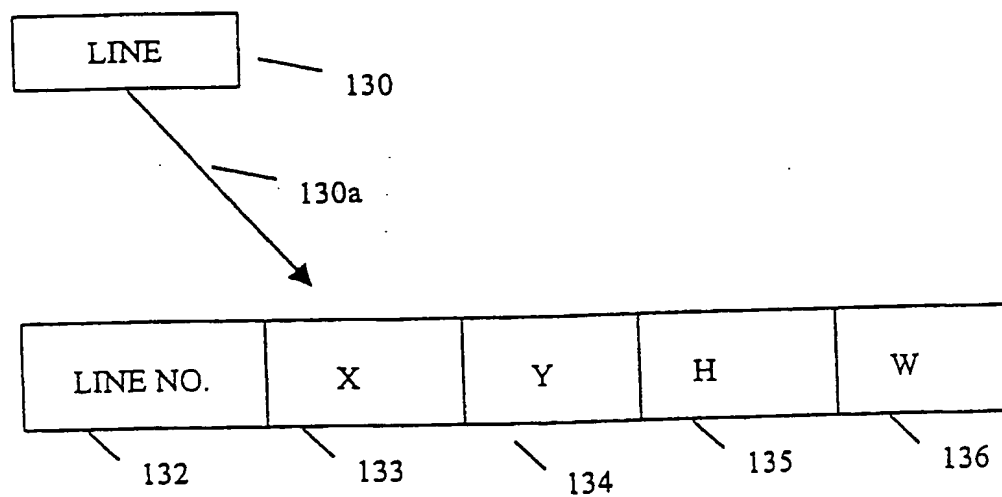


FIG. 8

13/17

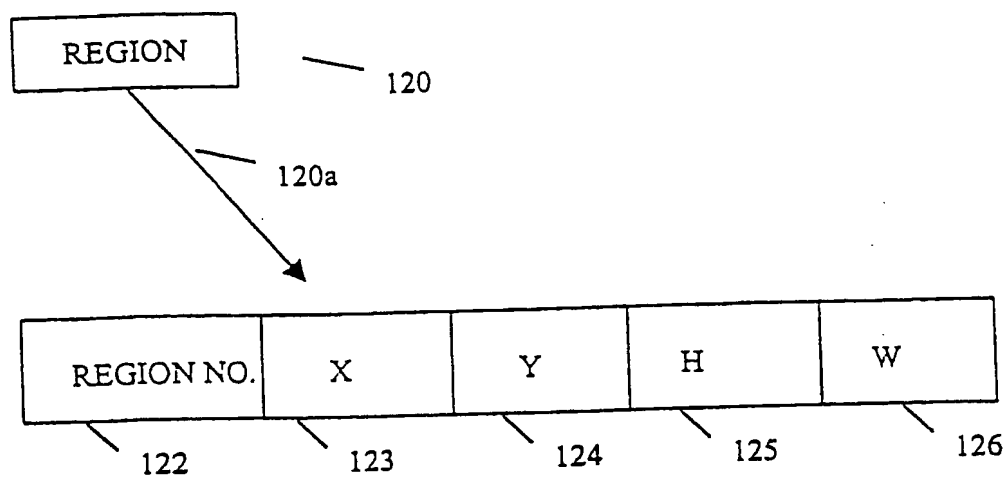


FIG. 9

14/17

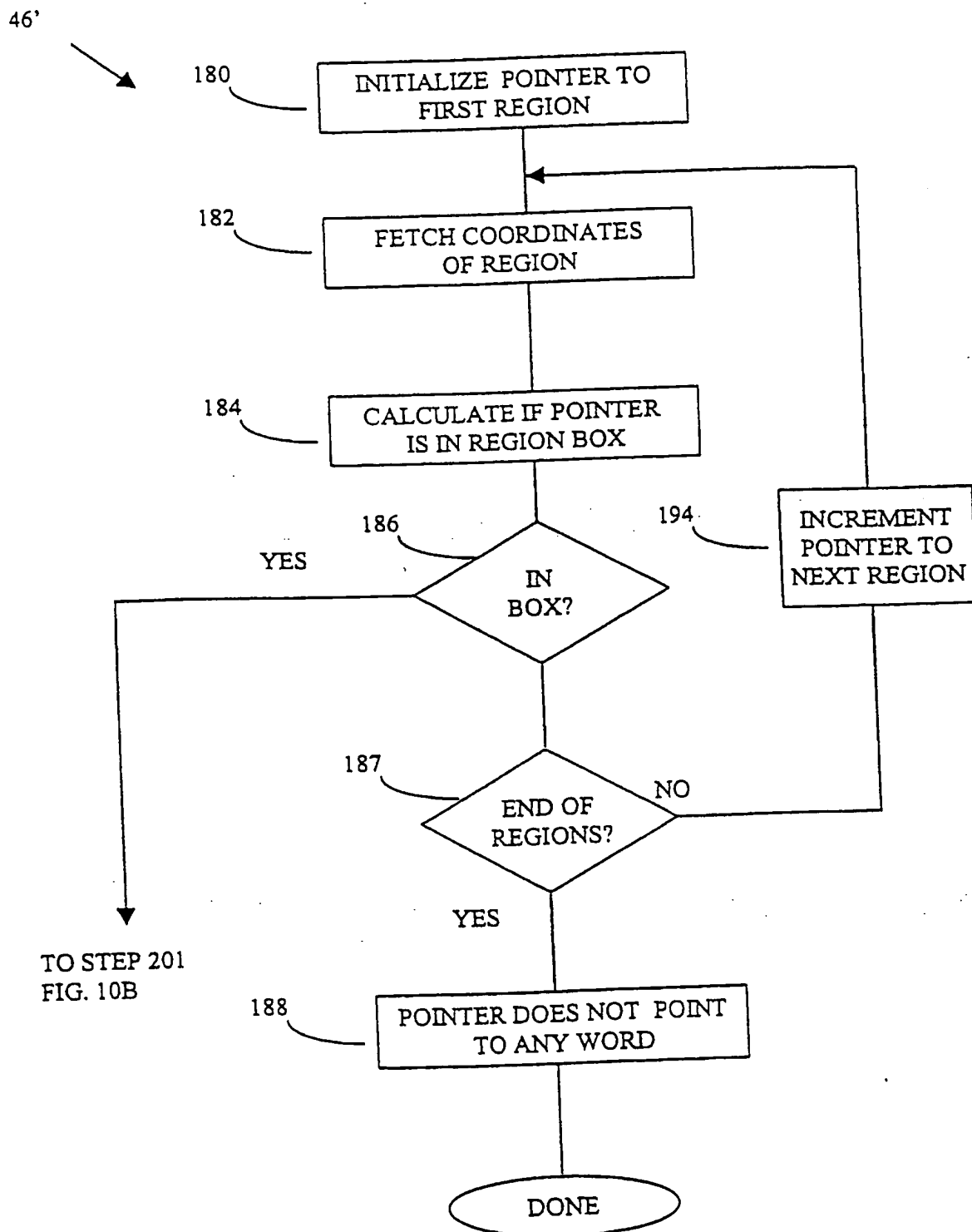
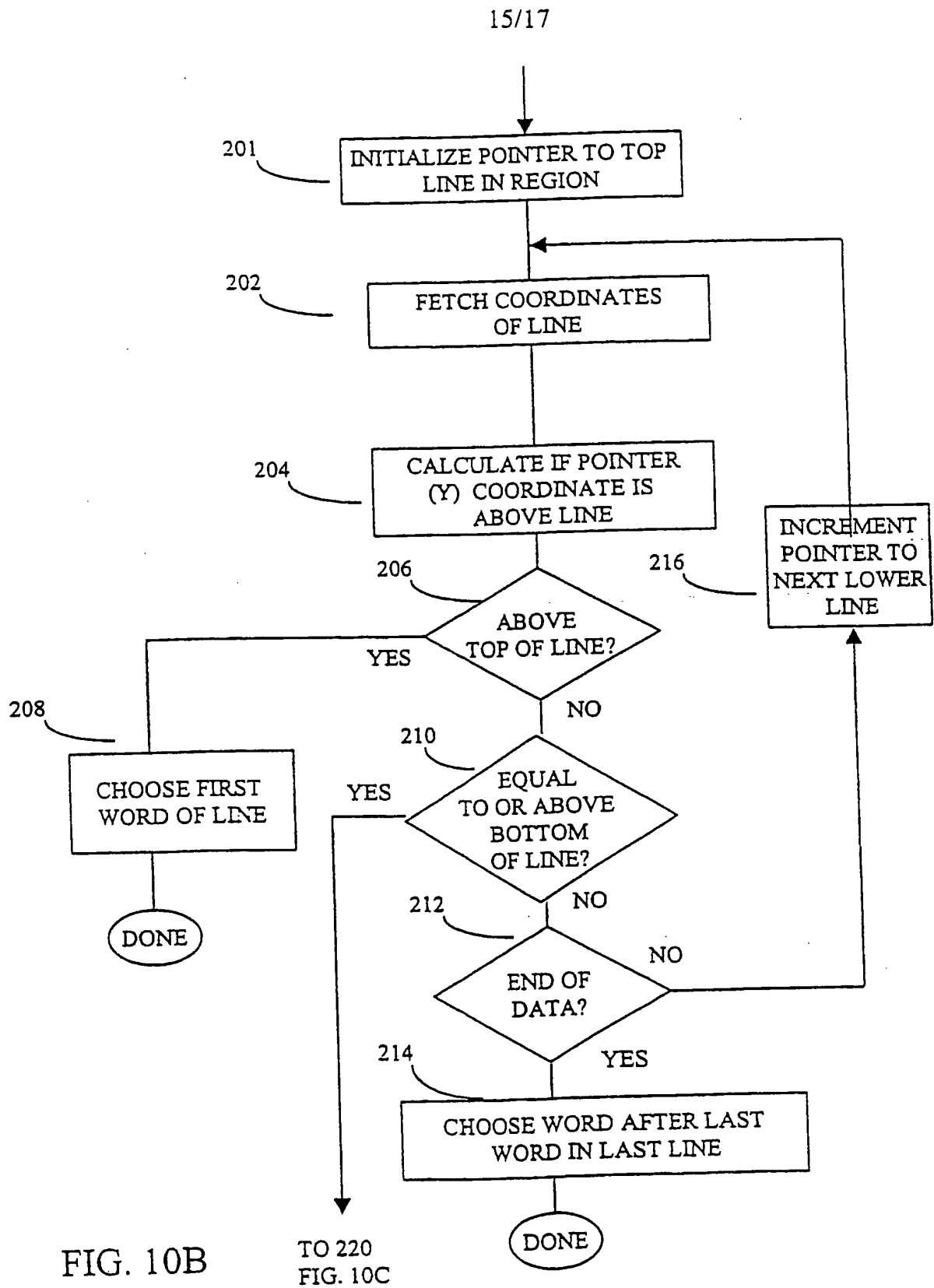


FIG. 10A





16/17

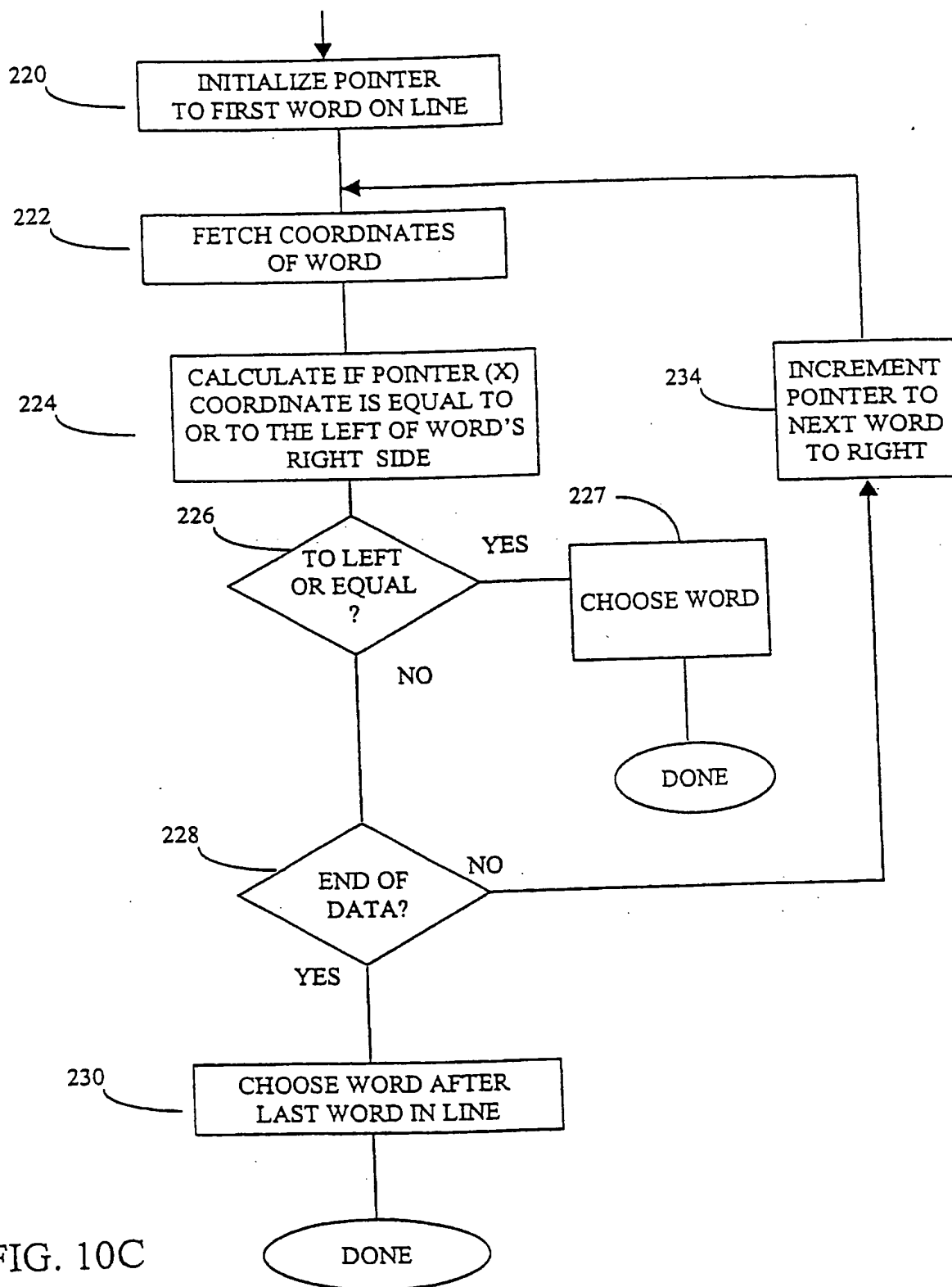


FIG. 10C

17/17

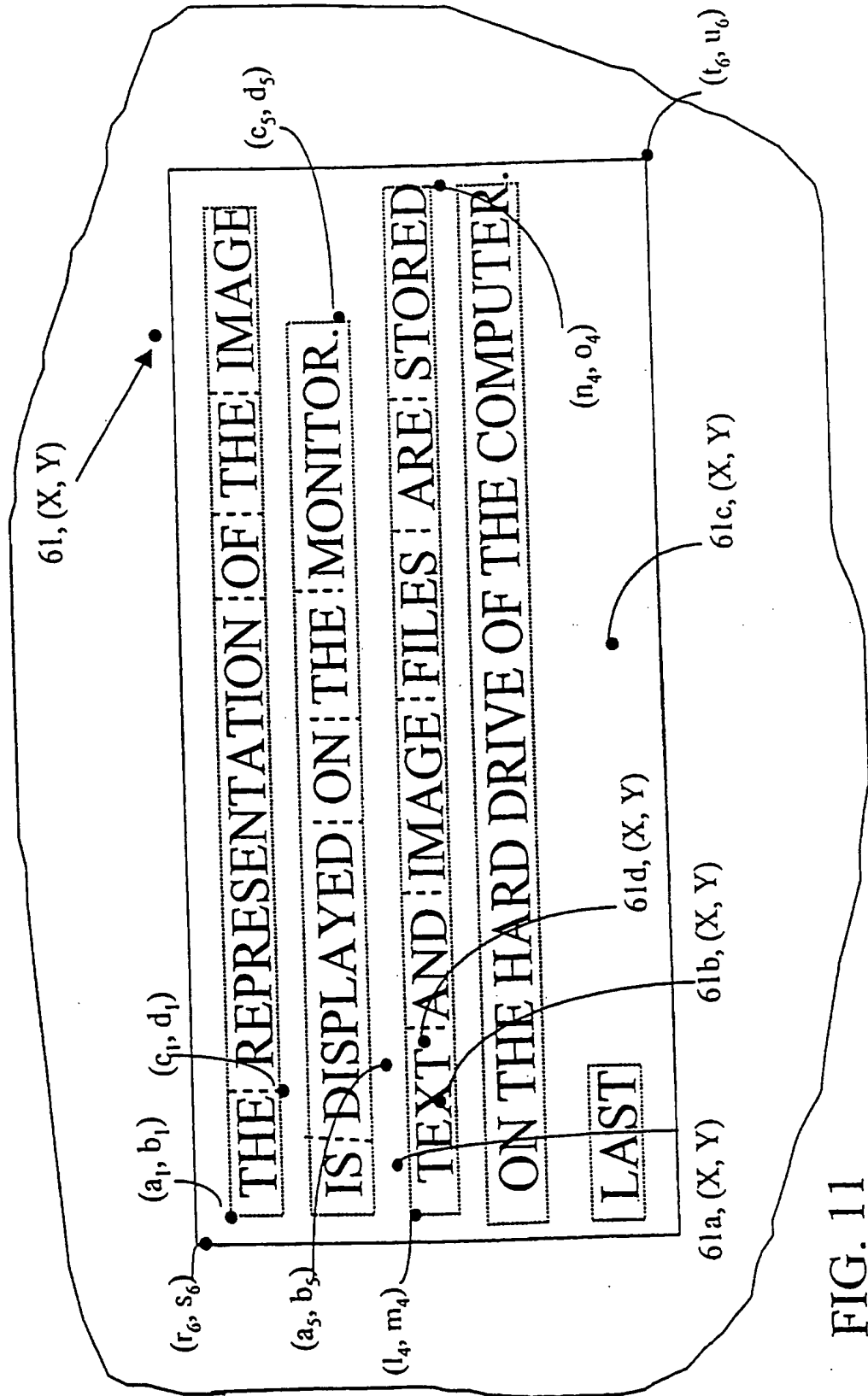


FIG. 11

## INTERNATIONAL SEARCH REPORT

International application No.

PCT/US99/13886

## A. CLASSIFICATION OF SUBJECT MATTER

IPC(6) : G10L 3/00; G09B 5/00

US CL : 704/260, 278, 434/317

According to International Patent Classification (IPC) or to both national classification and IPC

## B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

U.S. : 704/260, 258, 278, 243, 503, 504, 434/308, 309, 317, 318, 319

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)  
APS, DERWENT(WEST)

## C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X ----- Y	US 5,663,748 A [HUFFMAN et al] 02 September 1997, abstract, col. 13-16 and figure 3	1-5, 30-35  1-11, 19-36
Y	US 5,623,679 A [RIVETTE et al] 22 April 1997, col. 3-5 and 15	7-11, 25-29, 36
Y	US 4,698,625 A [MCCASKILL et al] 06 October 1987, col. 2 lines 5-33	20-29
Y	US 5,538,430 A [SMITH et al] 23 July 1996, abstract	11, 19-24
Y	US 5,617,507 A [LEE et al] 01 April 1997, abstract	24

☒ Further documents are listed in the continuation of Box C.
 ☐ See patent family annex.

* Special categories of cited documents	*T* later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
*A* document defining the general state of the art which is not considered to be of particular relevance	*X* document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
*E* earlier document published on or after the international filing date	*Y* document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
*L* document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)	*Z* document member of the same patent family
*O* document referring to an oral disclosure, use, exhibition or other means	
*P* document published prior to the international filing date but later than the priority date claimed	

Date of the actual completion of the international search 13 AUGUST 1999	Date of mailing of the international search report 21 OCT 1999
Name and mailing address of the ISA/US Commissioner of Patents and Trademarks Box PCT Washington, D.C. 20231 Facsimile No. (703) 305-3230	Authorized officer HAROLD ZINTEL Telephone No. (703) 305-2381

## INTERNATIONAL SEARCH REPORT

International application No.

PCT/US99/13886

## C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y	US 5,875,428 A [KURZWEIL et al] 23 February 1999, whole document	1-11,19-36